

Identifying Clusters using Node2Vec, Spectral Clustering and GCN

Aayush Shrivastava Vivek Tarachandani Ankit Pandey Anish Pawar Puneet Kumar Rajan
{cs24mtech02001, ai23mtech14008, ai23mtech13002, ai23mtech14002, ai22mtech11010}@iith.ac.in

Indian Institute of Technology, Hyderabad

Abstract

Node2Vec is a framework for learning continuous feature representations for nodes in networks. Spectral Clustering is a technique that uses eigenvalues of the similarity matrix of the data to reduce dimensions before clustering in fewer dimensions. This method is particularly effective for identifying clusters that are not necessarily globular and can capture complex cluster structures that traditional clustering methods might fail to detect. Graph Convolutional Networks (GCN) are powerful neural network architectures for processing data from graphs.

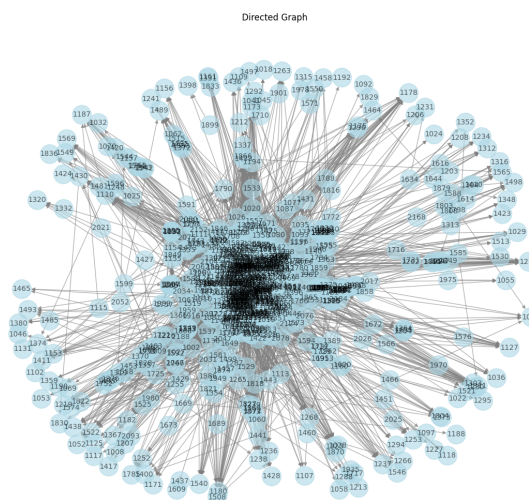
1. Problem Statement

The Node2Vec project aims to develop an algorithm that generates robust feature representations for network nodes using biased random walks, enhancing tasks like node classification and community detection while preserving network topology. The Spectral Clustering project aims to partition graphs into clusters based on the eigenvalues of their similarity matrices, enabling the identification of complex community structures that traditional clustering methods might miss. The GCN project seeks to develop a neural network architecture that leverages node features and graph structure to effectively perform tasks like node classification, link prediction, and graph classification within graph-based data. This approach aims to enhance the learning of node representations by incorporating neighborhood information directly into the learning process, enabling more accurate predictions and insights into the underlying graph structure.

2. Dataset Description

The dataset "Payments.csv" contains 1,30,535 rows and three columns: Sender (Integer), Receiver (Integer), and Amount (Integer). Each row represents a transaction where a certain amount is transferred from the Sender to the Receiver. Out of these transactions, only 5,358 involve a

unique pair of dealers interacting with each other. These 5,358 transactions are unique interactions between pairs of dealers who have conducted business with each other at least once.



represents the transaction between the sender and the receiver and the value of the transaction amount.

4. Node2Vec

Node2Vec framework is based on learning continuous feature representation for nodes in the graph and preserving the knowledge gained from the neighboring 100 nodes. Let us understand how the algorithm works. Say we have a graph with a few interconnected nodes creating a network. So, the Node2Vec algorithm learns a dense representation (say 100 dimensions/features) of every node in the network. The algorithm suggests that if we plot these 100 dimensions of each node in a 2-dimensional graph by applying PCA, the distance of the two nodes in that low-dimensional graph would be the same as their actual distance in the given network.

Algorithm 1: Node2VecWalk

Data: Graph $G' = (V, E, \pi)$, Start Node u , Length l
Result: A walk of length l starting from node u

```

1 Initialize walk to  $[u]$ 
2 for  $walk\_iter = 1$  to  $l$  do
3    $curr \leftarrow walk[-1]$ 
4    $V_{curr} \leftarrow \text{GetNeighbors}(curr, G')$ 
5    $s \leftarrow \text{AliasSample}(V_{curr}, \pi)$ 
6   Append  $s$  to walk
7 return walk

```

Algorithm 2: Learn_Features

Data: Graph $G = (V, E, W)$
Result: Feature representation f for all nodes in V

```

1 Parameters: Dimensions  $d$ , Walks per node  $r$ ,  

   Walk length  $l$ , Context size  $k$ , Return  $p$ , In-out  $q$ ;
2  $\pi \leftarrow \text{PreprocessModifiedWeights}(G, p, q)$ ;
3  $G' \leftarrow (V, E, \pi)$ ;
4 Initialize walks  $\leftarrow \emptyset$ ;
5 for  $iter \leftarrow 1$  to  $r$  do
6   forall  $u \in V$  do
7      $walk \leftarrow \text{Node2VecWalk}(G', u, l)$ ;
8     Append walk to walks;
9   end
10 end
11  $f \leftarrow \text{StochasticGradientDescent}(k, d, \text{walks})$ ;
12 return  $f$ ;

```

5. Spectral Clustering

Spectral Clustering is a variant of the clustering algorithm that uses the connectivity between the data points

to form the clustering. It uses eigenvalues and eigenvectors of the data matrix to forecast the data into lower dimensions space to cluster the data points. It is based on a graph representation of data where the data points are represented as nodes and an edge represents their similarity.

Algorithm 3: SpectralClustering(D, s, k)

Input: D : a set of data points $\{x_1, \dots, x_n\}$,
 $s(x, x')$: a similarity function, k : number of clusters to construct

Output: Clusters A_1, \dots, A_k

```

1 Construct a similarity graph  $G$  from  $D$  using  $s$ 
2 Let  $W$  be the weighted adjacency matrix of  $G$ 
3 Compute the graph Laplacian  $L$ 
4 Compute the first  $k$  eigenvectors  $u_1, \dots, u_k$  of  $L$ 
5 Let  $U \in \mathbb{R}^{n \times k}$  be the matrix containing  $u_1, \dots, u_k$ 
   as columns
6 for  $i = 1, \dots, n$  do
7   | Let  $z_i \in \mathbb{R}^k$  be the  $i^{th}$  row of  $U$ 
8 end
9 Cluster the points  $z_i$  with the k-means algorithm
   into clusters  $C_1, \dots, C_k$ 
10 for  $i = 1, \dots, k$  do
11   | Let  $A_i = \{x_j \mid z_j \in C_i\}$ 
12 end
13 return  $A_1, \dots, A_k$ 

```

6. Graph Convolutional Network (GCN)

Graph Convolutional Networks (GCNs) generalize convolutional neural networks to graph-structured data. They operate by learning a function that takes a graph's adjacency matrix, a degree matrix, and a feature matrix as input and produces node-level embeddings as output. These embeddings can then be used in various downstream tasks such as node classification, link prediction, or graph classification.

Algorithm 4: GraphConvolutionalNetwork (GCN)

Input: Graph $G = (V, E)$ with adjacency matrix A ,
degree matrix D , feature matrix X , depth
 K , non-linearity σ , learnable weight
matrices $W^{(k)}$, $k = 1, \dots, K$

Output: Node-level embeddings Z

```

1  $H^{(0)} \leftarrow X$ 
2 for  $k = 1, \dots, K$  do
3   |  $H^{(k)} \leftarrow \sigma \left( D^{-\frac{1}{2}} A D^{-\frac{1}{2}} H^{(k-1)} W^{(k)} \right)$ 
4 end
5  $Z \leftarrow H^{(K)}$ 
6 return  $Z$ 

```

7. Results

We obtain dense vector embeddings for each node upon applying the Node2Vec algorithm to our graph data. These embeddings capture the local neighborhood structure and encode the nodes' positions within the overall graph topology. Leveraging these embeddings, we utilize the k-means clustering algorithm to partition nodes into distinct clusters, ideally corresponding to different communities or roles within the graph.

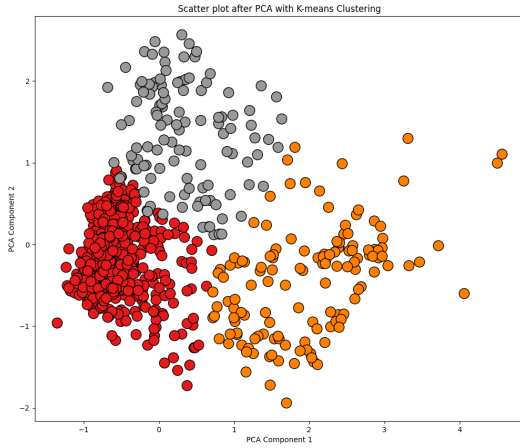


Figure 2: Node2Vec Output

In contrast, when employing Spectral Clustering, the approach capitalizes on the eigenvalues of the graph's Laplacian matrix to perform dimensionality reduction before clustering. This method is particularly adept at identifying non-linearly separable clusters in the original space, often revealing global structures and groupings that are not immediately apparent.

Graph Convolutional Networks (GCN), on the other hand, provide a more nuanced method of embedding generation. By iteratively aggregating and transforming neighbor information, GCNs produce embeddings that inherently incorporate both local graph topology and node features. These embeddings can then be clustered like those from Node2Vec to discover community structures.

8. Conclusion

In conclusion, the assignment employing Node2Vec, Spectral Clustering, and GCN techniques demonstrates the versatility and effectiveness of various graph-based approaches in uncovering latent structures within complex networks. Each method offers unique advantages: Node2Vec captures local neighborhood information, Spectral Clustering reveals global structures, and GCNs

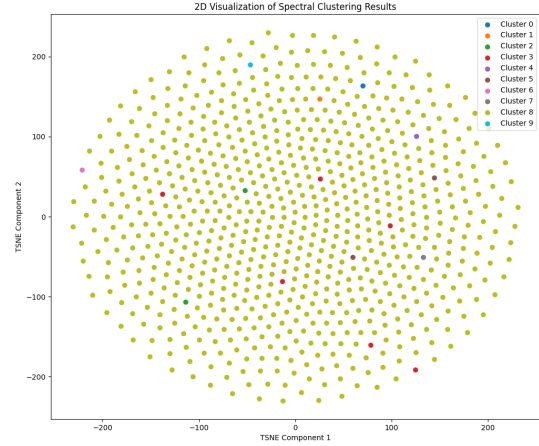


Figure 3: Spectral Clustering Output

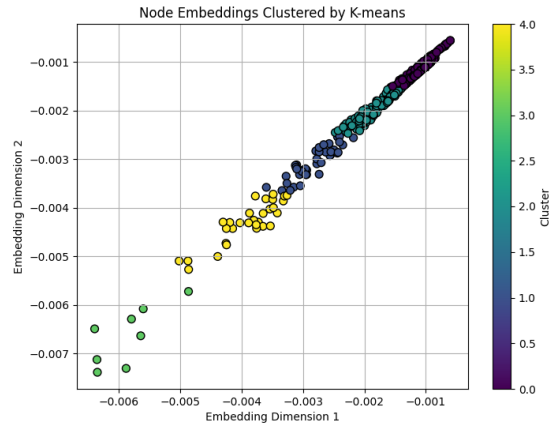


Figure 4: GCN Output

integrate node features and graph topology. The choice among these techniques depends on the specific characteristics of the data and the objectives of the analysis, showcasing the importance of tailored algorithm selection for optimal results in graph analytics tasks. According to the results, GCN provides the most accurate clustering, followed by Node2Vec and Spectral Clustering.

Node2Vec Code is available in the given link.
Spectral Clustering Code is available in the given link.
GCN Code is available in the given link.

References

- [1] Node2Vec: Scalable Feature Learning for Networks (2016), Aditya Grover, Jure Leskovec.

- [2] A Tutorial on Spectral Clustering (2007), Ulrike von Luxburg.
- [3] Semi-Supervised Classification with Graph Convolutional Networks (2016), Thomas N. Kipf.